An Architecture for Emotional Facial Expressions as Social Signals

Ruth Aylett[®], Christopher Ritter, Mei Yii Lim, Frank Broz, Peter E McKenna, Ingo Keller, and Gnanathusharan Rajendran

Abstract—We focus on affective architecture issues relating to the generation of expressive facial behaviour, critique approaches that treat expressive behaviour as only a mirror of internal state rather than as also a social signal and discuss the advantages of combining the two approaches. Using the FAtiMA architecture, we analyse the requirements for generating expressive behavior as social signals at both reactive and cognitive levels. We discuss how facial expressions can be generated in a dynamic fashion. We propose generic architectural mechanisms to meet these requirements based on an explicit mind-body loop and Theory of Mind (ToM) processing. A illustrative scenario is given.

Index Terms—Intelligent agents, affective computing, interactive systems, software architecture, cognitive informatics

12 **1** INTRODUCTION

1

5

6

7

g

10

11

13 [¬]HIS paper poses the problem of how to incorporate a generative account of expressive behaviour into an affective 14 architecture, focusing on facial expressions. Expressive behav-15 iour using the body posture, gesture, glance, facial expression 16 is an significant component of communicative content along-17 side the verbal channel, and is therefore required for social 18 19 agents, whether robots or graphical characters. Facial expressions are considered particularly important for agents that 20 have a face (some robots do not), since this is often the focus of 21 glance by an interaction partner. With more than forty muscle 22 groups [1], the face has a wide range of movements and thus 23 substantial expressivity. It has been argued that more than 24 half of expressive behaviour relates to facial expressions [2]. 25

In this paper we focus on facial expressions that relate to affect. Many computational accounts, when not using scripting, treat them as a mirror of the internal affective state of the agent and as a way of signalling that state to interaction partners [3]. This makes for an architectural mechanism that is conceptually straightforward: directly connecting the affective outputs of the architecture to the expressive modalities of

- R. Aylett is with MACS, Heriot-Watt University, Edinburgh EH14 4AS, United Kingdom. E-mail: r.s.aylett@hw.ac.uk.
- C. Ritter is with CITEC, Bieldefeld University, Bieldefeld 33615, Germany. E-mail: critter@techfak.uni-bielefeld.de.
- M.Y. Lim, F. Broz, and I. Keller are with the School of Mathematical and Computer Sciences, Heriot-Watt University, Edinburgh EH14 4AS, United Kingdom. E-mail: [M.Lim, f.broz, ijk1]@hw.ac.uk.
- P.E. McKenna and G. Rajendran are with the Department of Psychology, Heriot-Watt University, Edinburgh EH14 4AS, United Kingdom. E-mail: {p.mckenna, T.Rajendran}@hw.ac.uk.

Manuscript received 18 June 2018, revised 6 Mar. 2019, accepted 12 Mar. 2019, Date of publication 0 . 0000; date of current version 0 . 0000. (Corresponding author: Ruth Aylett.)

Recommended for acceptance by C. M. de Melo, C. Becker-Asano, D. C. Moffat, J. Parthemor, and D. D. Petters.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TAFFC.2019.2906200 the agent. However it is clear that even children aged 3-4 [4], 33 never mind adults, routinely modify their facial expressions 34 in a number of ways related to their social context. 35

There is an argument for emotional transparency in a 36 social agent. Cognitive appraisal theory suggests that emo-37 tions are generated when an event is appraised against a per-38 son's goals, with positive affect when events favour goals 39 and negative affect when they do not. A transparent display 40 may make a social agent's goals, and how far these are suc-41 cessfully met, more obvious to its interaction partner. Thus 42 the early system Kismet [5] played the role of an infant in 43 learning scenarios, and used transparent expressive behav-44 iour to regulate the interaction, to encourage more or reduce 45 the amount of stimulus it was receiving.

However many applications of embodied social agents 47 require more sophisticated expressive behaviour. The Laura 48 agent of [6] deliberately generated warm facial expressions 49 so as to build trust and rapport. Where an embodied social 50 agent aims to improve user motivation, or it operates in a 51 training or education setting [7] then expressive behaviour 52 is more likely to be an action explicitly chosen by the agent 53 than a reflection of its internal state. Issues relating to longlived interaction and getting past the novelty effect [8] also 55 require a less naive account of expressive behaviour. Dealing with these issues is a motivation for this work. 57

A simple case is where a facial expression related to internal affect is muffled or suppressed. A classic example would 59 be playing the card game poker, though a subordinate being 60 reprimanded by their boss, or a parent walking with their 61 child on a dark night, might be expected to suppress expressions of anger and fear respectively. In the case of poker, there 63 are game-related social norms; in the other two examples a 64 mix of social norm and in-situ estimation of the impact of 65 ones expressive behaviour on others in the shared context. 66

[9] distinguishes four categories of expression modification 67 (*display rules*): the cultural; the personal (depending on 68

personality or other individual factors); vocational require-69 ments (as in actors); and the need of the moment. Giving a 70 social agent this ability to suppress expressions requires a pro-71 cess tracking the contextual impact of a given expression. 72 Empathy, in which the observer responds to the affective state 73 of another is one means of carrying out this tracking. More 74 generally, one might posit Theory of Mind (ToM) capabilities, 75 in which one takes into account how one is likely to be per-76 ceived by another [10]. At the very least the agent must be able 77 to recognise and adhere to social norms for a given situation. 78

In the pedagogical example above, expressive behaviour 79 may also be generated for a specific communicative objective. 80 Likewise, economic games used to study negotiation may 81 deploy expressive behaviour competitively as game moves 82 [11] or as an aid to reaching cooperative decisions [12]. The 83 84 classic case of an unwelcome birthday present may result in a deliberate facial expression of pleasure in order to convey grat-85 86 itude. As cited above, children as young as 3 or 4 perform this modification. Indeed smiles are notoriously ambiguous about 87 internal state and very often related to social context [13]. A 88 recent account [14] gives three types of social smiles: those 89 90 rewarding the behaviour of others, as in this example, those creating or strengthening affiliative social bonds, and those 91 regulating social hierarchies. We should note that embodied 92 social agents without specific expressive behaviour may still 93 have their behaviour treated as socially expressive. [3] exam-94 ples a robot that turned away from a user immediately follow-95 ing a request and was interpreted as showing dislike or 96 contempt. Thus an embodied social agent that can predict the 97 social impact of its expressive behaviour may help to prevent 98 misunderstandings. We propose to do this through an explicit 99 mind-body loop and the application of ToM capabilities. 100

101 The key point is that facial expressions operate as social signals not merely as information about internal state. Even 102 103 the greater emotional transparency of infants relates to a social context where carers are motivated by smiles and will 104 act to deal with the causes of negative affective states. We 105 argue that coupling the affective outputs of an agents archi-106 tecture to its expressive modalities is insufficient in the devel-107 opment of embodied social agents. 108

By modelling the ability to handle expressive behaviour as 109 a social signal, we broaden the range of applications to which 110 a social agent can be applied and offer a standard mechanism 111 for regulating expressive behaviour - whether by suppression 112 or substitution - rather than ad hoc application-dependent 113 solutions. We also broaden its communicative repertoire by 114 explicitly including expressive behaviour in the set of actions 115 from which it can select rather than binding it to the internal 116 affective state being modelled. We seek to retain some of the 117 interactional advantages of transparent expressive behaviour 118 119 by supporting modification and not just overlaying of expressions, allowing micro-expressions [15] that we refer to else-120 where as the *Partial Poker Face* [16]. 121

However, should one equip social agents with what could 122 amount to deceptive behaviour? There are ethical considera-123 tions in creating convincing liars even though we know that 124 many humans behave like that [17]. But the line between 125 actually lying and what we describe as social facilitation is 126 very blurred. [18] argue that up to 30 percent of social interac-127 tions of longer than 10 minutes contain deceptions about 128 affective state. Thus more long-lived social agents do need 129

more sophisticated expressive behaviour to create smooth 130 interaction.

A second set of architectural requirements is raised by the dynamic nature of expressive behaviour. Some representations, such as the Facial Action Coding System (FACS) [1], define specific static facial expressions. In reality, facial expressions are nearly always continuous and dynamically varying, along with the speech stream they often accompany. Statebased expressive behaviour fits well with a state-based architecture, with dynamics confined to interpolating between the defined expressive states. In addition, a state-based approach fits well with explicit representations of affective state, where a dynamic approach to expressions is more consistent with a process-based architecture and implicit representations. We return to these issues in the next section.

2 BACKGROUND

2.1 Theoretical Issues

While we have argued above against wholly transparent 147 affective expressive behaviour, the idea that expressive 148 behaviour represents affective state in humans at all is far 149 from uncontested. [19] argues strongly the *behavioural eco-* 150 *logical view*, that facial expressions are not related to affec-151 tive state but are entirely social signals produced by an 152 evolutionary process. However we do not align with this 153 more radical viewpoint, being more convinced by argu-154 ments against it in [20] and studies such as [21] which 155 shows widespread interpretation of facial expressions as 156 indicative of affective state. 157

A more categorical position still is to reject the idea that 158 affective states cause actions at all, not just expressive 159 behaviour conceived as action. This relates to discussions of 160 Basic Emotion Theory (BET), that a finite set of emotions 161 such as anger, fear, disgust, happiness, sadness, surprise 162 and possibly others, emerge from evolutionary processes 163 related to survival and operate reflex behaviours. As [22] 164 argues in relation to his New BET, discussion is bedeviled 165 by using linguistic labels to mean different things, from pro-166 cesses at different levels of abstraction (e.g., physiological 167 sensations v cognitive categories) to different affective states 168 (are all forms of *anger* the same?). In this paper we take the 169 perspective of cognitive appraisal theory, that affective state 170 creates action tendencies priming actions rather than inevi- 171 tably producing them, while we also model lower-level pro- 172 cesses that have the character of reactions, if not reflexes. 173

A generic issue is how far an affective architecture can be 174 considered social rather than merely individual. Computa- 175 tional architectures based on psychological theory tend to 176 import individualist assumptions. The *Big Five* personality 177 dimensions [23] sometimes used for behaviour generation 178 in embodied social agents focuses on individual patterns of 179 behaviour. Cognitive appraisal theory [24] is not per se 180 incompatible with the modelling of social and cultural pro-181 cesses, but its focus on interaction between external events 182 (stimuli) and individual goals prioritises the individual. 183 Where the goals are taken as a given, affect will then represent an established individual reaction to the social context. 185

Indeed, [24] does not distinguish between appraisals relating to an event and to a person, so that socially-determined motions such as *sorry-for* (someone) are modelled in exactly 188

the same way as the *fear* generated by a threat to one's survival. Yet the social context is known to have a substantial
effect on individual appraisal: [25] gives the example of
watching a comedy you enjoy with a close friend who disapproves of it.

Social appraisal [25], [26] involves appraising the thoughts 194 195 or feelings of others, especially those with whom there is a relationship, as well as the emotion-causing event itself. This 196 view stresses the importance of empathic reactions and the 197 role of expressive behaviour in social regulation processes 198 which might result in modifications of expressive behaviour. 199 Note that social appraisal does not require an actual change 200 in internal affective state. Sensing the disapproval of a friend 201 - a social signal - one might actually find a comedy less amus-202 ing, or, for affiliative reasons, suppress the expression of one's 203 204 amusement.

Social appraisal is not unlike the idea of coping behav-205 206 iour [27]. Coping behaviour, a reaction to an affective state, has an external path in which actions-in-the-world are car-207 208 ried out to make the world more compatible with the goals of the individual. It also has an internal path, where cogni-209 tive strategies adjust a painful affective state, for example 210 by perceiving a 'silver lining' to an unpleasant event. Inter-211 nal coping behaviours are a second possible source of 212 expressive behaviour modification. 213

Cognitive appraisal architectures can be extended into a more socially responsive form by adding explicit mechanisms to handle social interaction. Thus the FAtiMA architecture [28] has been extended with a simulation ToM capability [29] (see Section 4), and the ability to model culturally-specific behaviour [30]. Both offer mechanisms supporting the modification of expressive behaviour.

221 The developmental robotics approach [31] is more likely to give the social context priority since it considers the con-222 223 struction of internal architectural structures by interactional processes such as enaction [32]. However since this work is 224 driven by the analysis of very young infants, most of that 225 looking at communication considers basic capabilities like 226 mutual glance and the development of turn-taking. For 227 facial expressions [33] the issue of interest is how to learn a 228 mapping between expressions and internal state. This does 229 not bear on the problem considered here. 230

A further theoretical dimension is a static versus a 231 dynamic account of behaviour, with implications for the 232 approach a computational implementation must take. In its 233 234 early form, cognitive appraisal was distinctly state-based: an event was compared with the goals of the individual, a set 235 (usually) of labelled emotions was generated to enable action 236 tendencies. More modern versions of cognitive appraisal, 237 such as the Component Process Model (CPM) [34], break 238 239 appraisal into a sequence of related actions in a pipeline of evaluation phases, each with a set of subchecks that may be 240 shared between phases. The main phases of the CPM concern 241 Relevance, Implication, Coping Potential, Normal Significance 242 243 where the last of these refers to social norms. Linking these to Facial Action Units, as CPM does in some cases, allows facial 244 expressions to be generated directly from the appraisal pro-245 cess [35] without having to pass through labelled emotions, 246 thus producing a more dynamic model. 247

A multi-stage model reduces the granularity of each step, moving in the direction of process. In addition, because facial

expressions can be driven at different stages on different 250 timescales, it supports micro-expressions and expression 251 modification. However multi-stage appraisal is not the only 252 way of dealing with this issue. A different class of models, 253 those based on drives and homeostasis, are more directly pro- 254 cess-based [36]. The PSI model [37] works with drives based 255 on five basic needs: personal survival (food etc); species sur- 256 vival (sex etc); affiliation (belonging to a group, engaging in 257 social interactions); certainty (the need for predictability of 258 events and consequences) and competence (the need to mas- 259 ter problems and tasks, including meeting needs). Drive- 260 based architectures work with upper and lower bounds on 261 needs, setting a comfort zone within which the values are 262 acceptable. If a need moves out of the comfort zone, the drive 263 seeks to activate behaviours to move it back - this is the pro- 264 cess of homeostasis which is inherently dynamic. 265

The PSI model has no direct representation of emotion, 266 but the behaviours generated by the drives are interpretable 267 as having affective qualities such as anger, joy or anxiety. It 268 outputs numerical values of valence and arousal which can 269 be used to synthesise multiple expressive behaviours with- 270 out having to pass through labelled emotions or necessarily 271 through facial action units [38]. The downside of a model at 272 this lower level of abstraction is that it is difficult to use for 273 embodied social agents using natural language, since 274 language by definition works at the symbolic level. This 275 suggests a multi-level model, very common in robotics, 276 in which moment-to-moment behaviour is controlled by 277 drives but strategic decisions are made via appraisal and 278 planning. This is the approach taken in the FAtiMA archi- 279 tecture discussed below. A multi-level model is also a multi-stage model, with lower layers typically working on 281 shorter timescales, thus also supporting micro-expressions 282 and expression modification. 283

2.2 Implementations

In this section we consider implemented systems that deal 285 with expression modification and also with facial expres- 286 sions generated as explicit communications. 287

One body of work on expression modification is concerned 288 with combining more than one emotion to produce mixed 289 facial expressions, for example an immediately generated 290 emotion with longer term affective states like mood [39]. This 291 is however still a version of transparency. [40] discusses the 292 issue of social modification of expressions but focuses on the 293 actual composition process. [41] takes this idea further by con- 294 sidering how different emotions might arise from an egocen- 295 tric appraisal and an empathic appraisal and evaluating the 296 impact on users of masking (empathic expression overrides 297 egocentric expression) and superimposition (expressions are 298 combined). However expressions were hand-coded rather 299 than generated autonomously by an architecture and the 300 focus was on realising and then evaluating the expressions in 301 a graphical character. 302

Empathic behaviour requires facial modification within a 303 social context. It is a significant issue in work on pedagogi- 304 cal agents, though some work [42], [43] relates to natural 305 language expressions rather than non-verbal behaviour. [7] 306 discusses robot expressive behaviour in a tutoring context, 307 but the robot used had no facial expressivity and relied on 308 gesture, while affective responses were derived from an 309

application-specific learner model that would not generaliseto other domains.

[44], in a role-play therapeutic domain, argues expressive 312 behaviour may relate to affect generated by an agent's own 313 coping actions, citing guilt as a result of a shift-blame action. 314 This produces a sequence of expressive behaviours but still 315 relates to the actual affective states of the agent. This appli-316 cation couples a cognitive architecture similar to the one we 317 modify to a pre-authored dialogue model; as it involves 318 agent-agent rather than human-agent dialogue, it can be 319 certain about the affective states each agent is responding 320 to. Thus it combines a cognitive architecture able to model a 321 rich internal state with the focus on interactivity of a dia-322 logue system, albeit a pre-authored one. 323

Interesting work on affect in interaction has been carried 324 325 out in the context of negotiation games. [11], [45] studies the social impact of an agent's display of joy, sadness, anger and 326 327 guilt, and how they function as social signals. These studies support learning of the parameters for a Bayesian network 328 giving probabilistic predictions, from emotion displays, of 329 how the negotiating partner appraises the interaction. This 330 supports predictions about their intentions [46]. This is not 331 intended to be and is not a generic architecture but is specific 332 to the iterated prisoner's dilemma used in the studies. 333

Other work on expressive behaviour as social signal 334 focuses on deception and lies. [17] investigates how deceptive 335 expressive behaviour may be used to produce a desired out-336 come in an economic game, using a similar approach to [45]. 337 Focusing on a particular element of negotiation, this work 338 demonstrates that agents with a deceptive facial expression 339 340 when they make an offer do attain the desired negotiation outcome. However the study was carried out with video record-341 342 ings of agent expressions rather than with a generating architecture. 343

344 [47] directly investigates expressive behaviour for an agent telling lies, building on [48]. It argues that facial expressions 345 will not be completely deceptive because some facial muscles 346 are not controllable at the conscious level. Thus both micro-347 expressions and compound facial expressions will occur. This 348 work also considers timing, with faked expressions lasting 349 longer than transparent ones and asymmetry, where faked 350 expressions have more activity on one side of the face than 351 the other. Two studies were carried out, which both used 352 smiles to mask other expressions in a similar way to [40]. In 353 the first study, smiles were either straightforwardly happy or 354 combined with disgust: these conditions were hand-coded. 355

A second study used a liar dice game in which lies are part of game play. As with other work using games, the context is easy to assess, depending entirely on the game play, so a generic architecture for deception was not required. In both cases, the aim was to evaluate the social impact of the compound expressions.

Work that is very relevant to this paper came from [49]. 362 This distinguished between an impulsive agent one that 363 364 directly expresses its affective state and a reflexive agent, which could refrain from expressing its state. However it 365 was based on the annotation of dialogue plans rather than 366 an affective architecture. These were held in a plan library 367 within the framework of a BDI architecture [50] originally 368 designed for rational agents and only later extended, in con-369 ceptual form, to incorporate affect [51]. 370

This view of expressive behaviour as a multi-modal adjunct 371 to language communication comes from a community with a 372 different focus from cognitive or affective modelling. It gener-373 alises the idea of a performative language action into non-374 verbal behaviour [52]. It has a strong focus on interactivity, but 375 the social signal aspects of expressive behaviour are inferred 376 from the dialogue moves with which it is associated. In this tra-377 dition, dialogue is viewed as a means of changing beliefs in a 378 purely logical model [53], an orthogonal view to cognitive 379 modelling. It supports affective communication but without a 380 modelled affective state or any affect-generating process.

Cognitive model-based research puts language actions on 382 a similar level to other agent actions rather than giving them 383 control of agent activity, while in dialogue system research, 384 agent actions are determined by a dialogue manager. This delgates the control of expressive behaviour to a process that 386 annotates utterances using a mark-up formalism ([54], [55]). 387

Mark-up of a dialogue stream both gives primacy to language over expressive behaviour, and assumes that the decisions about what to communicate are made by the Dialogue Manager before affective expressive behaviour is generated. However a social signal approach requires that affective choices be made at the level of action selection in the architecture. We also argue that the modification process requires an internal circuit within the architecture, since if the agent does not know what affective response it has chosen, it cannot modify it in a contextually sensitive manner. In human terms, you need to know you are angry in order to suppress anger.

This follows the work of [56] and in turn the ideas of [57] 399 who stresses the somatic aspects of emotion, an emotional 400 body state, which feeds into later processing at a more cognitive level. It is also consistent with the more sophisticated 402 view of cognitive appraisal already mentioned [58] as a multistage process with a temporal profile, consisting of cognitive 404 appraisal, a physiological activation and involving arousal, 405 motor expression (expressive behaviour), a motivational component, and a state of subjective feeling. If one sees expressive 407 behaviour as integral to emotion in this way, then even as a reflection of internal state it poses architectural issues.

Finally, [59] addresses some of the same questions as this 410 work but focuses on social signal analysis so as to establish 411 the affective state of an interaction partner rather than social 412 signal generation. It is concerned with recognising the social 413 signals associated with emotional regulation - or coping 414 behaviour - focusing on shame, generated in a cognitive 415 appraisal framework as a result of agent actions that have 416 negative praiseworthiness. It incorporates a simulation model 417 of ToM, similar to that discussed below in Section 4, imple-418 menting a shame model in its own architecture to make pre- 419 dictions about an interaction partner's likely social signals, 420 and thus aid Bayesian recognisers in detecting them. This 421 maintains a transparency approach. It does enrich its model 422 by including the target of the expressive behaviour - an issue 423 not yet dealt with in the work here. 424

3 REPRESENTATIONAL ISSUES

3.1 Architectures

All computational architectures are shaped by their 427 representational choices. We have already referred to one 428 significant dimension, the choice between state-based and 429

4

```
Fig. 1. An example of BML from SmartBody, controlling Glance.
```

process-based representations. In affective architectures,
this has tended to result in a choice between symbolic representations actuated through inferencing (for example [28],
[39] and non-symbolic representations in which homeostasis is a dominant mechanism (for example [37]).

Dialogue management systems were once symbolic sys-435 tems, manipulating natural language. More recently, after the 436 success of statistical approaches in speech recognition, 437 machine learning on large corpuses of spoken dialogue in spe-438 cific domains has encoded probabilistic transitions between 439 dialogue actions [60]. This makes the addition of expressive 440 behaviour, especially that related to affect, more difficult as it 441 requires analysis of the chosen dialogue action with fallible 442 approaches such as sentiment analysis. It produces an archi-443 444 tecture in which the role of inferencing is substituted by transitions in the learned network, an implicit encoding. Thus the 445 446 issue of explicit versus implicit representation is a further 447 dimension.

448 The most important representational decision in an affective architecture is how to represent affect itself, determined 449 in large part by the chosen theoretical framework. A simple 450 state-based architecture may represent affect as a single 451 symbolic variable and an associated intensity value, as in 452 the OCC model [24] with its 42 named emotions, while a 453 model built around drives and homeostasis may have no 454 explicit representation of affect at all [37] but generate affec-455 tive behaviour as dynamic patterns. A multi-stage theory 456 such as the one in this work may involve multiple represen-457 tations of affect, in particular if affect is seen as a phenome-458 non on more than one architectural level. 459

The FAtiMA architecture used as the basis for the ideas 460 below [28] divides into a reactive and a predictive compo-461 nent working on different time-scales and controlling differ-462 ent types of behaviour. Its predictive component runs a 463 464 planning system, which is where overt communicative intent would be handled. But some behaviour cannot rea-465 sonably be thought of as consciously planned - take the 466 example of bursting into tears at the death of a loved one. 467 468 FAtiMA incorporates a reactive layer that triggers much shorter-term unplanned behaviours. Note that incorporat-469 ing a reactive system does not in itself force the choice 470 between a symbolic or non-symbolic representation since 471 symbolic rules can play this role. 472

473 We have seen that cognitive appraisal itself may be 474 decomposed into multiple stages suggesting a collection of processes rather than a single process. There is physiologi- 475 cal evidence [61] of different brain mechanisms being 476 involved in what is known as emotional or affective empathy (or sometimes as emotional contagion), and cognitive 478 empathy, based on reasoning about the affective state of 479 another. 480

Damasio [57] distinguishes between primary and secondary emotions. The former are seen as innate, relating to fast and reactive behaviour patterns such as fight/flight, or infant distress behaviours, that do not involve cognitive-level procsend are closely tied to specific stimuli. Secondary emotions like *admiration* or *hope* are ascribed to higher cognitive processes involving expectations, learned outcomes and social context. Note that this distinction does not correspond exactly to the language labels that we may use: *fear* may count as a primary emotion if one is attacked by a ferocious dog, but we may use the same label in relation to an event that has not yet happened, an inverse to hope.

Primary emotions independent of the social context, are 493 good examples of emotions an agent might suppress after a 494 later evaluation. If a ferocious dog attacked while one was taking a child for a walk, a flight response activated by a primary 496 emotion of fear might be suppressed in favour of an attack on 497 the dog to defend the child. The language label for this would 498 be *courage*. However, the anger one might feel witnessing the 499 action of a bullying boss against a fellow employee is not primary by this definition, though its expression or suppression is also subject to evaluation of the social context. 502

While the distinction between primary and secondary 503 emotions is not wholly useful for this work, that between 504 somatic and cognitive impact does capture the stages to be 505 modelled so that an agent can *feel* an emotion so as to modify 506 it. The modelling issue is how to represent affect in cognitive 507 and somatic systems and how to link these representations 508 both to evaluation of social context and to the generation of 509 expressive behaviour. 510

3.2 Representations for Facial Expressions

Accounts of expressive behaviour that reduce it to multimodal annotation of the output from an affective architecture (or indeed, from a dialogue manager [62]) pose the representational problem as one of mapping from the architecture (conceived as mind) to the agents actuation capabilities (conceived as body). Interesting work in graphical characters has moved towards a standardised mark-up language for this purpose, Behavioural Markup Language (BML) [63] and to middleware based on this [64] such as SmartBody. Fig. 1 gives an example of BML controlling an agents gaze in Smart-Body [64].

The standard BML flow assumes that behaviour plan- 523 ning will deal with annotations on utterances from a 524 higher-level process. It leaves no space for a somatic repre- 525 sentation that can reflect emotion back into the cognitive 526 system to be re-evaluated. Such a causal chain would run as 527 shown in Fig. 2 producing a mind-body-mind loop. 528

The somatic level both dispatches output to the actual 529 generation of expressive behaviour and is also evaluated so 530 that the agent *knows what it is feeling* and is able to start the 531 modification process. The somatic level also supports 532 modelling of involuntary expressive changes - for example 533 blushing or crying. The temporal overlap produced by 534



Fig. 2. Affect re-evaluation.

re-evaluation is consistent with the idea of micro-expressions
sions [15] that facial expressions will reflect internal affective state for only a very short time before being replaced by
the socially determined expression. This idea is central to
our proposed architecture.

A somatic level also has a role to play in mapping affect onto an agent's expressive capabilities. Happiness can be expressed as a smile if a social agent can smile. If not - for example a robot with no face, or a face without a moveable mouth - other modalities can be selected.

Work has taken place to refine the markup system and to 545 add a specific virtual body representation [65] in the Thala-546 mus system. This inserted a BodyInterface unit between the 547 agent mind and behaviour planning in the same way as the 548 somatic level in Fig. 2. It has the ability to both receive and 549 send messages. This incorporates a feedback mechanism, 550 needed for the expressive capabilities under discussion, 551 but, as conceived, deals with external events and not inter-552 nal affective events 553

The two most widely-used systems for transforming an affective response into a behaviour specification in embodied agents are the Facial Action Coding System [1] and dimensional systems, of which the Pleasure-Arousal-Dominance System (PADS) [66] is the most popular.

FACS defines 44 muscle groups on the human face and 559 relates muscles to expressions via facial Action Units (AUs), 560 specific configurations of these muscle groups. It is often 561 used by researchers who want to generate facial expressions 562 in social agents, but is of much wider applicability. It was 563 designed for facial analysis, for example on videos, and is 564 still much more widely used for this purpose than for gener-565 ation. It is also used for research into expression recognition. 566

Used generatively, FACS offers a way of defining certain 567 facial expressions with respect to specific affective states, and 568 569 a tool is available for doing this [67]. In particular, AUs can be used to define Ekman's (contested) conception of primitive 570 emotions [68] facial expressions corresponding to fear, anger, 571 disgust, happiness, sadness, surprise that are said to be recog-572 nised across cultures (though this has been recently chal-573 574 lenged; see [69]). However, if these primitive emotions are seen as comparable to the primary emotions discussed above, 575 then they are targets for modification, and indeed they are 576 rarely visible in everyday adult social interaction. They thus 577 578 form a basis for blending, as in [47] discussed above.

An AU-defined facial expression representing an affective state is straightforward to interface to an affective architecture outputting such states. However, this is also a disadvantage, since it produces a static and rigid mapping. The very concept of an *expression*, as against a behaviour, works poorly in actual interaction as distinct from in a photograph or a video frame. AU definitions say nothing about how the face 585 moves into and out of an expression. 586

FACS can be used in a more dynamic manner. The decom-587 posed appraisal process of [58] discussed above associates 588 AU changes with its various stages. This has been implemented in some embodied social agents: [35] applied Sherer's 590 theoretical framework in a game environment. It used only 591 the labelled emotions Joy, Sadness, Guilt, Anger alongside 592 intermediate appraisal step changes, but there is no evalua-593 tion detailed. [70] used Hot Anger and Fear, but found many 594 questions about dynamics unanswered by the theory. [13] 595 investigated user perceptions of different smiles. It found 596 these were impacted by issues such as amplitude, duration, 597 onset and offset velocities and not just the AU. Recent work 598 in robot expression recognition by the authors also concluded 599 that the dynamics were very important in recognition [71]. 600

A final point is that even humanoid social agents usually 601 lack equivalents the full set of Action Units, with faces 602 much less expressive than a human's; even more true for 603 robots. Subtle affective behaviours van still be produced if 604 alternatives are sought - which may draw on film animation 605 as well as psychological theory. 606

Dimensional systems offer a numerical representation of 607 emotion in a space defined by two or more dimensional axes 608 but do not directly provide expressive behaviour. Indeed, 609 emotions can be represented as locations in a numerical Pleasure-Arousal-Dominance space (the PAD system), symbolic 611 labels attached to locations, and then used to drive AUs. A 612 more consistent approach to facial expressions in architectures using PAD would involve driving facial features (or 614 AUs) directly from the dimensional values. Here it is not necessary to translate to a symbolic affective label and then back 616 to numerically-driven motor action. 617

An example of work taking this approach, though using 618 Pleasure (Valence) and Arousal only, and not Dominance, is 619 that by Lim and Aylett [38]. This applied the drive-based PSI 620 architecture [37] in the context of a story-telling guide, and 621 directly linked the output valence and arousal values to a 622 simple 2D graphical face with limited expressive features. In 623 the absence of a single psychological theory linking valence 624 and arousal to specific facial features, this work adopted a 625 number of heuristics found in the literature affecting eyes, brows and mouth, thus bypassing linguistic labels for emo-627 tion. Fig. 3 shows part of the resulting facial feature space for 628 valence against high arousal (in this system high arousal was 629 630 0 and low was 1).

The advantage of driving expressive behaviour like this 631 within a process based architecture, is that behaviour will 632 be naturally dynamic, and the issues of merging different 633 expressive modalities are dealt with separately for each fea- 634 ture. The disadvantage is that it lacks the experimentally- 635 validated status of FACS.

If the somatic representation in use is not a symbolic one, 637 using PAD space to transform numerical triplets (P, A, D) 638 into a symbolic representation of affect allows the somatic 639 representation to be manipulated much more easily in cogni- 640 tive-level processing. This is useful since we will see that 641 social interaction theories are easier to represent symboli- 642 cally. A reactive system using drives can produce a rapid 643 affective output, using PAD space, feeding directly into 644 motor action, at the same time as outputting the nearest 645



Fig. 3. Facial expressions driven directly by valence and arousal.

symbolic label in the space to cognitive processing. Though
this is not how our illustrative architecture below works, we
point out in our conclusions that there may be very good reasons for taking this approach.

650 4 BASE ARCHITECTURE

Including a mind-body-mind loop for expressive behaviour
is independent of the detail of the architecture used. Any
architecture that represents mind and body components
and a structural inter-connection between them could take
up these ideas with different implementational details.

However the approach requires an architecture modelling
social interaction so that feelings returned to for re-evaluation
can be assessed in the social context. It also requires a ToM
mechanism able to assess expressed feelings for their impact on the interaction partner. There are few existing

architectures to choose from. Building a new architecture 661 from scratch is certainly possible, but in this paper the aim is 662 to explain clearly how to deal with expressive facial behav- 663 iour as a social signal, so using an existing architecture 664 reduces the size of the task. 665

For this reason, we start from the FatiMA architecture 666 already mentioned [28] a cognitive appraisal architecture 667 which has already been extended with social interaction functionalities [72], in particular the Social Importance Model of 669 Kemper [73] - see Fig. 4 and a simulation-based Theory of 670 Mind capability [29]. 671

FAtiMA deals with events along two time-scales: a reactive timescale without any intermediate processing, and a 673 deliberative timescale that allows for planning or other cognitive processing before selecting an action. These can be 675 seen in Fig. 4 to the right top and right bottom. The model 676 also includes a Memory component in which KB is a knowledge-base of the surrounding world and its objects, and AM 678 is an affective memory of past interactions supporting 679 mood modifications. The motivational state contains goals, 680 and activated goals/current intentions. 681

The reactive system is required for immediate expressive 682 behaviour unrelated to planning expressions of intense distress such as crying, or of involuntary laughing. Architectures that only implement affective transparency could deal 685 with the whole of expressive behaviour like this. However, 686 an advantage of using FAtiMA to discuss social signals is 687 that it also generates affective responses via its Deliberative 688 Layer supporting planned or other cognitively processed 689 expressive behaviour as well. 690

Neither layer in the existing architecture entirely captures 691 the issue under discussion of modifying expressive behav- 692 iour. We argue that modification can both act as a Delibera- 693 tive Layer activity and as a reaction to the agents internal 694 feeling of its affective state. 695

4.1 Social Importance Model

The Social Importance Model [72] is an example of integrat- 697 ing social rules into a cognitive appraisal model. Kemper [73] 698 focuses specifically on power and status, both of which are 699 contribute to the modification of expressive behaviour for 700 three of the categories categories cited by Ekman [9] culture, 701



Fig. 4. FAtiMA with Kemper modelling

vocational, and needs of the moment. We here summarise theimplementation and refer the reader to [72] for further detail.

In the implementated system of 4, Social Importance (SI) 704 represents an aggregated generalisation on the intuitive 705 meaning of status, since the SI an agent is willing to ascribe to 706 another may be influenced by inter-personal liking, group 707 membership, adherence to social norms, expertise, and per-708 sonal attributes such as wealth or strength. SI is not seen as a 709 static quantity but can be increased or decreased during the 710 course of social interactions. 711

The model contains three types of rules in its Reactive 712 713 Cultural Appraisal function: SI Attributions, SI Conferrals and SI Claims. An Attribution occurs when an agent meets 714 another agent and uses social rules to decide how much SI 715 it should have. Conferrals are associated with agent goals 716 717 and result in actions that acknowledge through behaviour the SI attribution an agent has given another. Expressive 718 719 behaviour is one example of a conferral mechanism: as in the example of looking pleased at a birthday gift even if the 720 agent does not like it. Conversely, an SI Claim is behaviour 721 carried out by the agent to assert its own SI in the eyes of 722 another agent, determinable using the Theory of Mind sys-723 tem discussed in the next section. 724

This architecture also includes a component for dealing 725 with items that are socially symbolic rather than merely func-726 tional. Examples include wearing specific clothing like even-727 ing dress, or presenting a bouquet of flowers to a soloist at the 728 end of a concert. Such items impact both on the agents moti-729 vations and its model of the motivations of others. This cre-730 ates extra inputs into Goal Selection in the Deliberative layer. 731 732 Finally the architecture can store specific plans relating to social rituals. These are defined as action sequences with a 733 734 specific social meaning that must be executed with a fixed order and content - for example greeting someone, whether 735 736 by shaking hands, bowing or kissing cheeks.

Relative SI has an obvious role in the modification of
expressive behaviour. If another agent has a very high level
of SI, then an agent is likely to suppress negative expressive
behaviour such as anger or distress. If two agents have
equivalent levels of SI as in two close friends then much less
modification of expressive behaviour is required.

Note that this SI model has its own input into expressive 743 behaviour as with other aspects of an agents internal pro-744 cesses. Social signals of disapproval or embarrassment 745 could be invoked by another claiming more SI than an agent 746 has attributed to them, while approval could be invoked by 747 another attributing the SI an agent has attributed to itself. In 748 these cases the agent generates a negative or positive affec-749 tive state but the extent to which this is expressed will 750 depend on the relative SIs involved 751

752 4.2 Theory of Mind

[10] defined Theory of Mind as the ability to infer the full 753 range of epistemic mental states of others, i.e., beliefs, 754 desires, intentions and knowledge. The abstractions we 755 make about the states of mind of others and consequently of 756 our own, is a mechanism that helps to make sense of their 757 behaviour in specific contexts and predict their next action. 758 A single-level theory of mind allows us to represent an 759 embodied agents beliefs about another embodied agents 760 beliefs and is the minimum needed to consider the impact 761



Fig. 5. Simulation of ToM through recursion.

of one's own expressive behaviour on someone else. Most 762 adults have at least a two-level ToM allowing them to think 763 about beliefs about another's beliefs about another's belief's 764 - who might well be you. 765

There are two conceptually different approaches to the 766 human theory of mind: the Theory-Theory approach (TT) 767 and the Simulation-Theory approach (ST). According to TT, 768 the mental state we attribute to others is not observable, but is 769 knowable through intuition and insight. In implementation, 770 this is achieved by using inference rules to reason about the 771 beliefs of others over an explicit model of the other. 772

On the other hand, ST claims that every person simulates 773 being another while trying to reason about their epistemic 774 state. This means that one can use the same structures and 775 processes used to update ones own beliefs and knowledge 776 to simulate those of others. In implementation, this involves 777 re-running the agent architecture as if for a different agent, 778 and this is conceptually straightforward for a cognitive 779 appraisal architecture such as FAtiMA [29]. 780

Let us assume that Agent1 (A1) is the agent carrying out 781 the ToM evaluation and Agent2 (A2) is the target of this 782 ToM. Then in general for an action X1 of A1 : 783

set X1 to be the event E1 for appraisal
 784

785

- 2) Run a copy of A1 on E1
- 3) Take X2 output by this new appraisal as the action of 786 A2 787

This recursive use of the agent architecture to simulate 788 ToM is shown in Fig. 5. 789

In this form A1 assumes A2 is exactly the same as them- 790 selves. However if one also applies the Social Importance 791 model just discussed with A1 and A2 interchanged, then 792 effectively the ToM is modified to take into account the differ- 793 ence in the social relationship from the point of view of A2. 794

This work was implemented [74] as part of a group-based 795 deception game, Werewolves, in which one agent is secretly 796 a were-wolf, able to kill other players in a segment where 797 everyone has their eyes closed. After the event, agents take it 798 in turns to accuse each other, and it is clear that the agent 799 playing the werewolf not only has to lie about their status but 800 accompany it with convincing expressive behaviour if it is to 801 play well. 802

5 UPDATED ARCHITECTURE

A number of updates to this architecture are needed in 804 order to add the capability of modifiable expressive behav- 805 iour. Here initial conceptual work has been carried out 806 under the banner of the Partial Poker Face [75] capturing 807



Fig. 6. FAtiMA modified to allow modified expressive behaviour.

the idea that expressive behaviour modification in humans
is rarely perfect. Fig. 6 shows the changes that would have
to be made to the FAtiMA architecture.

This is a slightly simplified version of 4 with some additions: Intrinsic Events, Virtual Body and Expressive Behaviour components, Partial Poker Face (PPF) and Expressive Behaviour (EB) Social Rules, Re-evaluation and ToM. Actions are planned sequences from the Deliberative Layer, Partial Poker Face and EB Social Rules are in fact part of the Reactive Layer, extracted here to make the diagram clearer.

In order to motivate these changes, we work through a sce-818 nario from [72] used to discuss the SI model of Section 4.1. In 819 820 this scenario, a traveller enters a bar after failing to find the way to their hotel. At the start of the scene, there are only two 821 822 characters sitting in the bar and they are talking to each other. 823 The barman is absent (although he later appears). The goal of 824 the traveller is to find directions to their hotel. In the version 825 discussed in [72], the traveller is an avatar directed by a human user, and the discussion focused on the behaviour of 826 the two agents in the bar. We adapt it to investigate how the 827 expressive behavior of a social agent version of the traveller 828 would be generated in the proposed architecture. 829

830 5.1 Intrinsic and Extrinsic Events

The first change that is needed to construct the mind-body-831 mind circuit at issue is a distinction between events external 832 to the agent entering into cognitive appraisal (extrinsic) and 833 834 events within the agent (intrinsic) generated by its affective responses. This distinction was not made in the ToM discus-835 sion above because the original motivation for that work was 836 allowing an agent to evaluate the impact on others of its exter-837 nal actions upon the world. Expressive behaviour resulting 838 839 from affective responses in agents exhibiting transparency had so far been considered to be hard-wired. As social signals 840 they can be thought of as responses to intrinsic events. Note 841 that intrinsic events may not be entirely invisible to other 842 843 agents where they are associated with truly involuntary behaviour at the physiological level, such as blushing or 844 sweating. 845

In the scenario, the traveller asks the bar agents if they
can give directions to the traveller's hotel. In the discussion
of [72], if these agents come from a collectivist culture, they
will be offended by the request since the traveller is not in

their in-group, and they will scowl and tell the traveller to 850 wait for the bar man to arrive. 851

In our simulation, the traveller agent will then appraise ⁸⁵² the rejection of their request as a goal failure, and within the ⁸⁵³ reactive layer interpret the scowls as disapproval. This will ⁸⁵⁴ generate a negative emotion of anger which is dispatched to ⁸⁵⁵ the Virtual Body. This corresponds to the somatic component of 2. The associated expressive behaviour is sent out to ⁸⁵⁷ the expressive system but the feeling of anger from the Virthe expressive system but the feeling of anger from the Virstual Body is passed back to the Intrinsic Event (IE) component, tagged with the Extrinsic Event (EE) ID and timestamp that originated it. This raises an intrinsic event, with ⁸⁶¹ an ID and time-stamp, making the agent aware of its emotion, thus beginning the mind-body-mind loop. ⁸⁶³

5.2 Re-Evaluation and Partial Poker Face

Conceptually, modifiable expressive behaviour must operate on a rapid reactive level as well as on a slower deliberative level. Otherwise modification would occur quite late 867 and the underlying emotion would result in clearly identifiable expressive behaviour. The speed of this reactive layer 869 is a variable in the personal presentation of an embodied 870 agent - some agents might suppress initial facial expressions 871 much faster than others, just as the intensity of the initial 872 emotion might vary between individual agents too and be 873 therefore harder or easier to suppress. 874

Thus the event signalling the traveller's anger is fed back 875 through the re-evaluation module to the PPF component 876 which is a set of reactive rules about dealing with emotions. 877 These rules draw on the agent's knowledge of social expressive behaviour which are flattened versions of ToM reasoning - that is to say compile the output of ToM reasoning into simple rules. In this case, PPF draws on a social rule that says if an agent expresses anger to someone it causes an angry response with high undesirability. Because this is a reactive rule it deals with a generic social situation, where the ToM module would reason about the concrete circumstances.

This rule leads PPF to genrate a neutral expression action 886 to suppress the angry expressive behaviour already dis-887 patched. The PPF sends this to the Action module where it 888 returns to the Virtual Body and is dispatched as expressive 889 behaviour. When it returns to Instrinsic Event but its IE tag 890 shows that it has been dealt with, preventing an infinite loop. 891

892 5.3 Deliberative Level

The IE component has also passed the emotion to the ToM module which is able to reason about the actual impact of an angry expression on the current goal of the traveller. ToM runs a copy of the appraisal system to assess the impact of the traveller's angry emotion on the bar agents, taking the SI of the traveller and the bar agents into account.

The ToM system will assess the angry expression as reducing the liking of the the bar agents for the traveller, reducing its SI and increasing the threat to its goal of finding the hotel.

This information is passed back to the deliberative system 902 whose planner assesses the ask-the-barman subgoal as a way 903 of achieving the find-hotel goal and deals with the SI threat 904 by proposing a smile and agreement with the proposal of the 905 bar agents. These actions are dispatched to the virtual body 906 907 but as the smile represents a social signal rather than an affective change it is not passed back round the mind-body-mind 908 909 loop by the Expressive Behaviour Component but only dispatched for execution. 910

The net effect on the traveller agent's behaviour is thus an angry micro-expression, suppressed quickly by a neutral expression and then replaced by a social smile.

914 6 CONCLUSIONS

In this paper we have tried to reformulate the architecture 915 of an embodied social agent to support the capability of 916 expressive behaviour as a social signal rather than only as 917 an indicator of the agent's internal state. The first question 918 one should ask is of course whether this is worth doing. 919 920 There are after all valid arguments for using an embodied agents expressive behaviour to reveal internal state, espe-921 922 cially in the case of robots, which are heavy pieces of metallic machinery that should share a human social space in a 923 924 way that is comfortable for humans. Knowing how a robot 925 is responding and what it is about to do are useful aspects of human-robot interaction. Indeed there is work [76] sug-926 gesting that action-expression revealing the motivation and 927 context for an agent action through expressive behaviour is 928 a necessity for smooth interaction. 929

As with so many questions in research, whether this is 930 worth doing is a case of *it depends*. Primarily it depends on 931 the social context in which the embodied agent is expected 932 to perform. We have already seen above that there are appli-933 cations for which modified expressive behaviour is not only 934 935 interesting but essential. This seems particularly true for those in which an embodied social agent is following a 936 vocational role, such as tutor, trainer, guide, receptionist. It 937 would be even more true if the application domain related 938 to drama and not to more naturalistic interaction, not only 939 940 in entertainment applications but also in areas such as roleplay based education and training. Here the activity itself is 941 limited by the complexity and expense of supplying human 942 actors, and there is clear scope for the use of social agents. 943

The approach discussed also supports in a principled way expressive behaviour which is difficult to generate - as against hardcode - without it. One group of such behaviours involves a combination of physiological signals with more cognitively-generated behaviour. Embarrassment, signalled by blushing (a physiological reaction) plus glancing away, would be an example of this. The blush can be generated very rapidly by the intrinsic event raised by the simulated 951 body, while the glance-away is generated later by consciously 952 'feeling' the emotion as it progresses further through the 953 mind-body-mind loop. A second group of behaviours relate 954 to the overlay of one expression by another as a socially deter-955 mined expression fails to completely override an internally 956 generated emotion. This supports the known issues with 957 smiles, which often combine with elements of other facial 958 expressions, such as the *disgust* hardcoded in by [40]. This is 959 achieved by a slow decay on a high-intensity emotion dis-960 patched from the simulated body and an overlaid smile from 961 the cognitive stage of the mind-body-mind loop. 962

Much of the discussion above - very much in line with the 963 literature - has taken a naturalistic approach, using normal 964 human social behaviour as a yardstick. However this 965 assumption should on occasion be challenged. It is not a fore-966 gone conclusion that this is the way to incorporate an embod-967 ied social agent into everyday human environments. These 968 are agents, they do not and will not for the foreseeable future 969 have human-level abilities given the extreme difficulties 970 involved. It could be that drama rather than naturalism is the 971 more useful paradigm. Indeed the idea of action-expression 972 [76] is more closely related to drama than to naturalism. By 973 showing a sequence of expressions as expressive behaviour 974 is modified, one supplies the human interaction partner with 975 information about the social adjustment of the agent, much in 976 the style of drama, where double-takes and slow realisations 977 are very much standard tropes. 978

A further argument in favour of a machinery for modify- 979 ing expressive behaviour is its use in decoding the expressive 980 behaviour of human interaction partners. The problems asso-981 ciated with facial expression recognition have not been the 982 subject of this discussion, but one of the most significant is 983 moving from sensor-based detection of facial movements to 984 an identification of the social signal being deployed. An agent 985 that has no concept of facial expressions as social signals, and 986 works on the basis that there will be a one-to-one mapping 987 between the expression and the users affective state is 988 unlikely to be successful. As argued at the start of this paper, 989 one can recognise a smile, but the signal the smile represents 990 is a different matter. An agent that has a simulation ToM 991 implementation can at least run its own architecture as a 992 decoding mechanism in the current social context. 993

6.1 Limitations

The most obvious limitation in the discussion of this paper is 995 that it is conceptual and has not yet been implemented. However, though much of the necessary basis for an implementation already exists in the FAtiMA architecture, the main point 998 being made here is that research into embodied social agents 999 should move from a widespread view of agent expressive 1000 behaviour as transparently affective especially in the case of 1001 facial expressions - and move to the social signal paradigm. 1002 We would argue that this also means a move away from the 1003 individualistic assumptions underlying many agent architectures into a more socially located account. 1005

A limitation of the suggested changes above in the 1006 FAtiMA architecture is that this version of the architecture 1007 is entirely symbolic in representation, making truly dynamic 1008 expressive behaviour problematic. In order to implement 1009 the dynamic PAD-space control of expressive behaviour 1010

10

discussed in Section 3 one would have to choose the FAtiMA variant FAtiMA-PSI [[77] system discussed above in Section 2.2. Here the FAtiMA symbolically-encoded reactive system is replaced by the PSI [37] five drives: Energy, Integrity, Affiliation, Certainty and Competence and a homeostatic mechanism that chooses actions and goals according to which drives need to be returned to their comfort-zone.

As an example of a non-symbolic approach that has no 1018 1019 explicit representation of affect, it is easy to see how it can drive expressive behaviour dynamically. It is less easy to 1020 see how one could incorporate the reactive elements of the 1021 Social Importance Model. While we argue that most of the 1022 discussion is relatively independent of the actual implemen-1023 tation architecture, it is clear that this difficulty would apply 1024 to other non-symbolic architectures, such as those generated 1025 1026 by machine learning approaches.

The recent ALEXA challenge [78], involving a disembod-1027 1028 ied (voice-only) conversational agent indicates the most 1029 likely solution. The systems that did best in an unrestricted 1030 conversational context were compound ones in which machine-learning derived transition networks sat under-1031 1032 neath symbolic rule-based systems that provided context and a degree of sanity check. It seems plausible that a com-1033 pound of this type could supply fast dynamic facial expres-1034 sions from its sub-symbolic processing, use PAD space to 1035 translate these into symbolic representations that are then 1036 passed into symbolic SI rules, and pass the outcome back 1037 through PAD space into the non-symbolic system. 1038

1039 In conclusion, these are the generic requirements for 1040 expressive behaviour as social signals outlined here.

- A mind-body-mind loop that allows an agent to feel its affect and can trigger.
- 1043 2) Intrinsic events differentiated from extrinsic events
 1044 coming from the surrounding environment.
- A re-evaluation process that responds to intrinsic
 events and modifies expressive behaviour
- A model of social interaction that can be used by the
 re-evaluation system to translate from desired social
 signal to modified expressive behaviour
- 10505)A ToM that can assess the social impact of an agents1051expressive behaviour and support a deliberative1052processing level in modification.

1053 We hope that this paper will help to stimulate work in 1054 improving expressive behaviour and changing the default 1055 approach to one of social signal generation.

1056 **ACKNOWLEDGMENTS**

The authors would like to thank members of the project ECUTE (ICT-5-4.2257666), whose work was partially supported by the European Commission (EC). The authors are solely responsible for the content of this publication. It does not represent the opinion of the EC, which is not responsible for any use that might be made of data appearing therein.

1063 **REFERENCES**

- 1064 [1] P. Ekman and W. V. Friesen, "The Facial Action Coding System," Qd5 Consulting, 1978.
- 1066 [2] A. Mehrabian, Nonverbal Communication. Evanston, IL, USA: Routledge, 2017.

- [3] M. F. Jung, "Affective Grounding in Human-Robot Interaction," in 1068 Proc. ACM/IEEE Int. Conf. Human-Robot Interaction, 2017, pp. 263–273. 1069
- P. M. Cole, "Children's spontaneous control of facial expression," 1070 Child Develop., vol. 57, pp. 1309–1321, 1986. 1071
- [5] C. Breazeal, "Role of expressive behaviour for robots that learn 1072 from people," *Philosophical Trans. Roy. Soc. B: Biological Sci.*, 1073 vol. 364, pp. 3527–3538, 2009.
- [6] T. Bickmore, "Relational agents: Effecting change through humancomputer relationships," *SciencesNew York*, 2003.
 1076
- [7] L. Hall, C. Hume, S. Tazzyman, A. Deshmukh, S. Janarthanam, 1077 H. Hastie, R. Aylett, G. Castellano, F. Papadopoulos, A. Jones, 1078 L. J. Corrigan, A. Paiva, P. A. Oliveira, T. Ribeiro, W. Barendregt, 1079 S. Serholt, and A. Kappas, "Map reading with an empathic robot 1080 tutor," in *Proc. 11th ACM/IEEE Int Conf. Human-Robot Interaction*, 1081 2016, pp. 567–567. 1082
- [8] R. Gockley, A. Bruce, J. Forlizzi, M. Michalowski, A. Mundell, 1083
 S. Rosenthal, B. Sellner, R. Simmons, K. Snipes, A. C. Schultz, and 1084
 J. Wang, "Designing robots for long-term social interaction," in 1085
 Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst., 2005, pp. 1338–1343. 1086
- [9] P. Ekman and W. V. Friesen, "Unmasking the face," Ann. Phys., 1087 1975.
- [10] D. Premack and G. Woodruff, "Does the chimpanzee have a theory of mind?" *Behavioral Brain Sci.*, 1978.
- [11] C. M. D. Melo, P. Carnevale, and J. Gratch, "The Effect of Expres- 1091 sion of Anger and Happiness in Computer Agents on Negotia- 1092 tions with Humans," in *Proc. 10th Int. Conf. Auton. Agents* 1093 *Multiagent Syst.*, 2011, pp. 937–944.
- [12] R. Hoegen, G. Stratou, and J. Gratch, "Incorporating emotion per- 1095 ception into opponent modeling for social dilemmas," in *Proc.* 1096 16th Conf. Auton. Agents MultiAgent Syst., 2017, pp. 801–809.
 1097
- [13] Z. Ambadar, J. F. Cohn, and L. I. Reed, "All smiles are not created 1098 equal: Morphology and timing of smiles perceived as amused, 1099 polite, and embarrassed/nervous," *J. Nonverbal Behavior*, vol. 33, 1100 pp. 17–34, 2009.
 [14] P. M. Niedenthal, M. Mermillod, M. Maringer, and U. Hess, "The 1102
- [14] P. M. Niedenthal, M. Mermillod, M. Maringer, and U. Hess, "The 1102 Simulation of Smiles (SIMS) model: Embodied simulation and the 1103 meaning of facial expression," *Behavioral Brain Sci.*, vol. 33, 1104 pp. 417–433, 2010. 1105
- [15] P. Ekman and W. V. Friesen, Unmasking the Face: A Guide to Recognizing Emotions from Facial Clues. Ishk, 2003. 1102
- [16] C. Ritter and R. Aylett, "The partial poker-face," in *Proc. Int. Conf.* 1108 *Intell. Virtual Agents*, 2015, pp. 479–482. 1109
 [17] J. Gratch, Z. Nazari, and E. Johnson, "The Misrepresentation 1110
- [17] J. Gratch, Z. Nazari, and E. Johnson, "The Misrepresentation 1110 Game: How to win at negotiation while seeming like a nice guy," 1111 in Proc. Int. Conf. Auton. Agents Multiagent Syst., 2016, pp. 728–737. 1112
- [18] B. M. DePaulo, S. E. Kirkendol, D. A. Kashy, M. M. Wyer, and 1113
 J. A. Epstein, "Lying in everyday life," J. Personality Soc. Psychol. 1114
 ogy, vol. 70, pp. 979–995, 1996. 1115
- [19] A. J. Fridlund, Human Facial Expression. New York, NY, USA: 1116 Academic, 1994. 1117
- [20] P. Ekman, "Should we call it expression or communication?" Inno- 1118 vation: Eur. J. Soc. Sci. Res., 1997. 1119
- [21] G. Horstmann, "What do facial expressions convey: Feeling 1120 states, behavioral intentions, or action requests?" *Emotion*, vol. 3, 1121 pp. 150–166, 2003. 1122
- [22] A. Scarantino, "Do emotions cause actions, and if so how?" *Emo-* 1123 *tion Rev.*, 2017. 1124
- [23] L. R. Goldberg, "An Alternative "Description of personality": The 1125 big-five factor structure," J. Personality Soc. Psychology, vol. 59, 1126 pp. 1216–1229, 1990. 1127
- [24] A. Ortony, G. L. Clore, and A. Collins, *The Cognitive Structure of* 1128 *Emotiions*. Cambridge, U.K.: Cambridge Univ. Press, 1988.
 1129
- [25] A. S. Manstead and A. H. Fischer, "Social appraisal: The social 1130 world as object of and influence on appraisal processes," in 1131 *Appraisal Processes in Emotion: Theory, Methods, Research*. London, 1132 U.K.: Oxford Univ. Press, 2001.
- B. Parkinson, A. H. Fischer, and A. S. Manstead, *Emotion in Social* 1134 *Relations: Cultural, Group, and Interpersonal Processes*. Evanston, IL, 1135 USA: Routledge, 2004.
- [27] R. S. Lazarus and S. Folkman, "Transactional theory and research 1137 on emotions and coping," *Eur. J. Personality*, 1987.
 1138
- [28] J. Dias and A. Paiva, "Feeling and reasoning: A computational 1139 model for emotional characters," in *Proc. Portuguese Conf. Artif.* 1140 *Intell.*, 2005, pp. 127–140. 1141
- [29] J. Dias, R. Aylett, A. Paiva, and H. Reis, "The great deceivers: Vir-1142 tual agents and believable lies," in *Proc. Annu. Meeting Cognitive* 1143 *Sci. Soc.*, 2013, vol. 35, no. 35.

- 1145 [30] S. Mascarenhas, J. Dias, R. Prada, and A. Paiva, "A dimensional 1146 model for cultural behavior in virtual agents," Appl. Artif. Intell., 1147 vol. 24, no. 6, pp. 552–574, 2010.
- [31] M. Asada, K. Hosoda, Y. Kuniyoshi, H. Ishiguro, T. Inui, 1148 1149 Y. Yoshikawa, M. Ogino, and C. Yoshida, "Cognitive Develop-1150 mental Robotics: A Survey," IEEE Trans. Auton. Mental Develop., vol. 1, no. 1, pp. 12-34, May 2009. 1151
- D. Vernon, "Enaction as a conceptual framework for developmental 1152 [32] cognitive robotics," Paladyn, J. Behavioral Robot., vol. 1, pp. 89-98, 1153 1154 2010
- 1155 [33] A. Watanabe, M. Ogino, and M. Asada, "Mapping facial expres-1156 sion to internal states based on intuitive parenting," J. Robot. 1157 Mechatronics, vol. 19, 2007
- [34] K. R. Scherer, "Appraisal considered as a process of multilevel 1158 sequential checking," 2001. 1159
- M. Courgeon, C. Clavel, and J.-C. Martin, "Appraising emotional 1160 [35] events during a real-time interactive game," in Proc. Int. Workshop 1161 Affective-Aware Virtual Agents Soc. Robots, 2009, Art. no. 7. 1162
- 1163 [36] L. Canamero, "Issues in the design of emotional agents," Cybern. 1164 Syst., 2001.
 - [37] D. Dörner and C. D. Güss, "PSI: A computational architecture of cognition, motivation, and emotion," Rev. General Psychology, vol. 17, no. 3, 2013, Art. no. 297.
 - M. Y. Lim and R. Aylett, "An emergent emotion model for an [38] affective mobile guide with attitude," Appl. Artif. Intell., vol. 23, no. 9, pp. 835-854, 2009.
- 1171 S. Marsella and J. Gratch, EMA: A Computational Model of Appraisal [39] 1172 Dunamics. 2006
- [40] M. Ochs, R. Niewiadomski, C. Pelachaud, and D. Sadek, 1173 "Intelligent expressions of emotions," in Proc. Int. Conf. Affective 1174 1175 Comput. Intell. Interaction, 2005, pp. 707-714.
 - [41] R. Niewiadomski, M. Ochs, and C. Pelachaud, "Expressions of empathy in ECAs," in Proc. 8th Int. Conf. Intell. Virtual Agents, 2008, pp. 37-44.
 - [42] S. W. McQuiggan and J. C. Lester, "Modeling and evaluating empathy in embodied companion agents," Int. J. Human Comput. Stud., vol. 65, pp. 348-360, 2007
- J. Robison, S. McQuiggan, and J. Lester, "Evaluating the conse-1182 [43] 1183 quences of affective feedback in intelligent tutoring systems," in Proc. 3rd Int. Conf. Affective Comput. Intell. Interaction Workshops, 1184 1185 2009, pp. 1-6.
- [44] S. C. Marsella, W. L. Johnson, and C. M. Labore, "Interactive Peda-1186 1187 gogical Drama for Health Interventions," in Proc. 11th Int. Conf. Artif. Intell. Educ., 2003. 1188
- 1189 [45] C. M. De Melo, P. Carnevale, and J. Gratch, "The effect of virtual 1190 agents' emotion displays and appraisals on people's decision 1191 making in negotiation," in Proc. Int. Conf. Intell. Virtual Agents, 2012, pp. 53-66. 1192
- 1193 [46] C. M. D. Melo, P. Carnevale, S. J. Read, and J. Gratch, "Bayesian 1194 model of the social effects of emotion in decision-making in multiagent systems," in Proc. 11th Int. Conf. Auton. Agents Multiagent 1195 Syst., 2012, pp. 55-62. 1196
- M. Rehm and E. André, "Catch me if you can exploring lying 1197 [47] agents in social settings," Aamas, 2005. 1198
- P. Ekman, "Darwin, Deception, and Facial Expression," Ann. New 1199 1200 York Acad. Sci., vol. 1000, pp. 205-221, 2003.
- [49] B. De Carolis, C. Pelachaud, I. Poggi, and F. De Rosis, "Behavior 1201 planning for a reflexive agent," in Proc. 17th Int. Joint Conf. Artif. 1202 Intell.- Vol. 2, 2001, pp. 1059-1064. 1203
- [50] A. Rao and M. Georgeff, "Modeling rational agents within a BDI-1204 1205 architecture," in Readings in Agents. San Francisco, CA, USA: Mor-1206 gan Kaufmann, 1997.
- [51] D. Pereira, E. Oliveira, N. Moreira, and L. Sarmento, "Towards an 1207 architecture for emotional BDI agents," in Proc. Portuguese Conf. 1208 1209 Artif. Intell., 2005, pp. 40-46.
- 1210 [52] I. Poggi and C. Pelachaud, "Emotional meaning and expression in animated faces," in Proc. Int. Workshop Affective Interactions, 2000, 1211 pp. 182–195. 1212
- C. Castelfranchi and Y. H. Tan, "The role of trust and deception in 1213 [53] 1214 virtual societies," Int. J. Electron. Commerce, 2002.
- [54] B. DeCarolis, C. Pelachaud, I. Poggi, and M. Steedman, "APML, a 1215 mark-up language for believable behavior generation," Lifelike 1216 Characters Tools Affective Functions Appl., 2004. 1217
- M. Schröder, "The SEMAINE API: Towards a standards-based 1218 [55] framework for building emotion-oriented systems," Adv. Human-1219 1220 Comput. Interaction, vol. 2010, 2010, Art. no. 2.

- [56] C. Becker-Asano and H. Ishiguro, "Evaluating facial displays of 1221 emotion for the android robot Geminoid F," in Proc. Workshop 1222 Affective Comput. Intell., 2011, pp. 1-8. 1223
- [57] A. R. Damasio, Descartes' Error. New York, NY, USA: Random 1224 House, 2006. 1225
- [58] K. R. Scherer, "On the nature and function of emotion: A component process approach," Approaches Emotion, vol. 2293, 1984, Art. no. 317. 1228
- [59] P. Gebhard, T. Schneeberger, T. Baur, and E. André, "MARSSI: 1229 Model of appraisal, regulation, and social signal interpretation," 1230 in Proc. 17th Int. Conf. Auton. Agents MultiAgent Syst., 2018, 1231 pp. 497-506. 1232
- [60] O. Lemon and O. Pietquin, "Machine Learning for Spoken Dia-1233 logue Systems," in Proc. INTERSPEECH, 2007. 1234
- [61] S. G. Shamay-Tsoory, J. Aharon-Peretz, and D. Perry, "Two sys-1235 tems for empathy: A double dissociation between emotional and 1236 cognitive empathy in inferior frontal gyrus versus ventromedial 1237 prefrontal lesions," Brain, vol. 132, pp. 617-627, 2009. 1238
- [62] O. Lemon, A. Bracy, A. Gruenstein, and S. Peters, "The WITAS 1239 multi-modal dialogue system I," in Proc. 7th Eur. Conf. Speech Com-1240 mun. Technol., 2001. 1241
- S. Kopp, B. Krenn, S. Marsella, A. N. Marshall, C. Pelachaud, [63] 1242 H. Pirker, K. R. Thórisson, and H. Vilhjálmsson, "Towards a com-1243 mon framework for multimodal generation: The behavior markup 1244 language," in Proc. Int. Workshop Intell. Virtual Agents, 2006, 1245 pp. 205–217. 1246
- M. Thiebaux, S. Marsella, A. N. Marshall, and M. Kallmann, [64] 1247 "SmartBody: Behavior realization for embodied conversational 1248 agents," in Proc. 7th Int. Joint Conf. Auton. Agents Multiagent Syst, 1249 2008, pp. 151–158. 1250
- [65] T. Ribeiro, M. Vala, and A. Paiva, "Thalamus: Closing the mind-1251 body loop in interactive embodied characters," in Proc. Int. Conf. 1252 Intell. Virtual Agents, 2012, pp. 189-195. 1253
- A. Mehrabian, "Analysis of the big-five personality factors in 1254 [66] terms of the PAD temperament model," Australian J. Psychology, 1255 vol. 48, no. 2, pp. 86–92, 1996. 1256
- [67]E. B. Roesch, L. Tamarit, L. Reveret, D. Grandjean, D. Sander, and 1257 K. R. Scherer, "FACSGen: A tool to synthesize emotional facial 1258 expressions through systematic manipulation of facial action uni-1259 ts,^{*/} J. Nonverbal Behavior, vol. 35, pp. 1–16, 2011. [68] P. Ekman and W. V. Friesen, "Constants across cultures in the 1260
- 1261 face and emotion," J. Personality Soc. Psychology, vol. 17, no. 2, 1262 pp. 124-129, 1971. 1263
- [69] R. E. Jack, O. G. B. Garrod, H. Yu, R. Caldara, and P. G. Schyns, 1264 "Facial expressions of emotion are not culturally universal," Proc. 1265 Nat. Acad. Sci., vol. 109, pp. 7241-7244, 2012. 1266
- [70] L. Malatesta, A. Raouzaiou, K. Karpouzis, and S. Kollias, "Tow-1267 ards modeling embodied conversational agent character profiles 1268 using appraisal theory predictions in expression synthesis," Appl. 1269 Intell., vol. 30, pp. 58-64, 2009. 1270
- [71] P. E. McKenna, M. Y. Lim, A. Ghosh, R. Aylett, F. Broz, and 1271 G. Rajendran, "Do you think I approve of that? Designing facial 1272 expressions for a robot," in Proc. Int. Conf. Soc. Robot., 2017, 1273 pp. 188–197. 1274
- S. Mascarenhas, N. Degens, A. Paiva, R. Prada, G. J. Hofstede, [72] 1275 A. Beulens, and R. Aylett, "Modeling culture in intelligent virtual 1276 agents," Auton. Agents Multi-Agent Syst., vol. 30, no. 5, pp. 931-962, 1277 2016. 1278
- [73] T. D. Kemper, Status, Power and Ritual Interaction: A Relational 1279 Reading of Durkheim, Goffman and Collins. Evanston, IL, USA: Rout-1280 ledge, 2016. 1281
- [74] R. Aylett, L. Hall, S. Tazzyman, B. Endrass, E. André, C. Ritter, 1282 A. Nazir, A. Paiva, G. Höfstede, and A. Kappas, "Werewolves, 1283 cheats, and cultural sensitivity," in Proc. Int. Conf. Auton. Agents 1284 Multi-Agent Syst., 2014, pp. 1085–1092. C. Ritter and R. Aylett, "The Partial Poker-Face," in Proc. Int. Conf. 1285
- [75] 1286 Intell. Virtual Agents, 2015, pp. 479-482. 1287
- [76] P. Sengers, "Do the thing right: An architecture for action-1288 expression," in Proc. 2nd Int. Conf. Auton. Agents, 1998, pp. 24-31. 1289
- M. Y. Lim, J. Dias, R. Aylett, and A. Paiva, "Creating adaptive 1290 [77] affective autonomous NPCs," Auton. Agents Multi-Agent Syst., 1291 vol. 24, no. 2, pp. 287-311, 2012. 1292
- [78] A. Ram, R. Prasad, C. Khatri, A. Venkatesh, R. Gabriel, Q. Liu, 1293 Nunn, B. Hedayatnia, M. Cheng, A. Nagar, and others, 1294 "Conversational AI: The science behind the Alexa prize," 1295 arXiv:1801.03604, 2018. 1296

1165

1166

1167

1168

1169

1170

1176 1177

1178

1179

1180

1181

AYLETT ET AL.: AN ARCHITECTURE FOR EMOTIONAL FACIAL EXPRESSIONS AS SOCIAL SIGNALS





Ruth Aylett became a professor of Computer Science at Heriot-Watt University in 2004 where she is a member of the Edinburgh Centre for Robotics and researches social agents, human-robot interaction and affective computing. She currently leads a UKRC-funded project SoCoRo investigated the development of a robot trainer in social signal recognition for high-functioning adults with an ASD.



Peter E McKenna received the doctorate degree1326in psychology from Heriot-Watt University, Edin-1327burgh. He is currently a research associate on1328the EPSRC funded SoCoRo project - he and the1320team are developing a socially competent robot1330to teach adults with an ASD social and employ-1331ment skills.1322



Christopher Ritter received the German diploma in computer sciences from Friedrich-Alexander University Erlangen-Nuremberg. He is currently working toward the doctoral degree in computer sciences at the University of Bielefeld, Social Cognitive Systems group, CITEC.



Ingo Keller received the diploma in computer sci-1333ence from Technische Universitt Dresden (TUD),1344Germany, in 2010. He joined Heriot-Watt University1335sity in 2014 to pursue the PhD degree and is1336researching interactive object learning in the area1337of Teachable Robots. He is also investigating1338aspects of gesture synthesis in the SoCoRo1339project.1340



Mei Yii Lim received the PhD degree from Heriot-Watt University, in 2007. She has worked as a post-doc researcher on EU-funded projects eCIR-CUS, eCUTE, SOCIETIES, LIREC and EMOTE. She is currently a researcher on the UK-funded project SoCoRo.



Gnanathusharan Rajendran received the PhD 1341 degree in developmental psychology from the 1342 University of Nottingham. He is an associate professor in psychology with Heriot-Watt University. 1344 He joined Heriot-Watt University in 2012 as a 1345 reader and specialises in typical and atypical cognitive and social development, digital education, 1347 and social robotics. 1348



Frank Broz received the PhD degree in robotics from the Carnegie-Mellon University Robotics Institute. He is an assistant professor of computer science with Heriot-Watt University. He was a senior research fellow with Plymouth University working on the Robot-ERA project before joining Heriot-Watt University in 2015 and researches artificial intelligence, human-robot interaction and social robotics. ▷ For more information on this or any other computing topic, 1349 please visit our Digital Library at www.computer.org/publications/dlib. 1350